

Claims

- 1 1. A method of cooperatively load-balancing a cluster of server computer  
2 systems for servicing client requests issued with respect to a plurality of client  
3 computer systems, said method comprising the steps of:  
4 a) selecting, by a client computer system, a target server computer system  
5 from said cluster of server computer systems to service a particular client request  
6 using available accumulated selection basis data;  
7 b) evaluating, by said target server computer system, said particular client  
8 request to responsively provide instance selection basis data dynamically  
9 dependent on the configuration of said target server computer and said particular  
10 client request; and  
11 c) incorporating said instance selection basis data into said available  
12 accumulated selection basis data to affect the subsequent selection of said target  
13 computer system with respect to a subsequent instance of said particular client  
14 request.
- 1 2. The method of Claim 1 wherein said instance selection basis data includes  
2 a representation of a dynamically determined performance level of said target  
3 server computer system and wherein said available accumulated selection basis  
4 data incorporates said instance selection basis data with identifications of said  
5 target server computer and said particular client request.
- 1 3. The method of Claim 2 wherein said instance selection basis data includes

2 a representation of a policy evaluation of said particular client request relative to  
3 said target server computer system.

1 4. The method of Claim 1 wherein said instance selection basis data includes  
2 a load value and a selection weighting value, wherein said load value represents  
3 a dynamically determined performance level of said target server computer system  
4 and said selection weighting value represents a policy evaluation of said particular  
5 client request relative to said target server computer system and wherein said  
6 available accumulated selection basis data incorporates said instance selection  
7 basis data with identifications of said target server computer and said particular  
8 client request.

1 5. The method of Claim 4 wherein said step of selecting selects said target  
2 server computer system based on predetermined selection criteria including the  
3 relative values of said load value and said selection weighting value with respect  
4 to said particular client request as recorded in said available accumulated  
5 selection basis data.

1 6. The method of Claim 5 wherein said instance selection basis data provides  
2 for a rejection of said particular client request and wherein said step of selecting  
3 includes selecting an alternate server computer system from said cluster of server  
4 computer systems as said target server system to service said particular client  
5 request based on said available accumulated selection basis data.

1 7. A method of load-balancing a cluster of server computer systems in the  
2 cooperative providing of a network service, said method comprising the steps of:  
3 a) selecting, by each of a plurality of host computers, server computers  
4 within a computer cluster to which to issue respective service requests;  
5 b) responding, by a corresponding one of said plurality of host computers,  
6 to the rejection of a predetermined service request by selecting a different server  
7 computer to which to issue said predetermined service request;  
8 c) receiving, in regard to said respective service requests by the respective  
9 ones of said plurality of host computers, load and weight information from the  
10 respective server computers; and  
11 d) evaluating, by each of said plurality of host computers, the respective  
12 load and weight information received with respect to server computers of said  
13 computer cluster as a basis for a subsequent performance of said step of  
14 selecting.

1 8. The method of Claim 7 further comprising the step of determining said  
2 weight information by each of said server computers with respect to each service  
3 request received, said weight information being determined from a predefined  
4 policy association between a received service request and the identity of the one  
5 of said server computers that receives the service request.

1 9. The method of Claim 8 further comprising the step of distributing initial  
2 information by said cluster of server computers to said host computers, said initial  
3 information providing selection lists of said server computers to said host  
4 computers.

1 10. The method of Claim 9 wherein said load information is representative of  
2 a plurality of load factors including network loading and processor loading.

1 11. The method of Claim 10 wherein said load information is representative  
2 of the processing of a current set of service requests including a plurality of  
3 processor functions.

1 12. The method of Claim 11 wherein said load information includes one or  
2 more load values representing processing functions internal to a server computer.

1 13. A server cluster operated to provide a load-balanced network service, said  
2 server cluster comprising:

3 a) a plurality of server computers individually responsive to service requests  
4 to perform corresponding processing services, wherein said server computers are  
5 operative to initially respond to said service requests to provide load and weight  
6 values, wherein said load and weight values represent the current operating load  
7 a policy-based priority level of a respective server computer relative to a particular  
8 service request; and

9 b) a host computer system operative to autonomously issue said service  
10 requests respectively to said plurality of server computers, said host computer  
11 system further operative to select a target server computer from said plurality of  
12 server computers to receive an instance of said particular service request based  
13 on said load and weight values.

1 14. The server cluster of Claim 13 wherein said host computer is operative to  
2 collect said load and weight values from said plurality of server computers in  
3 connection with the issuance of respective service requests to said plurality of  
4 server computers and wherein the selection of said target server computer is  
5 based on the relative temporal age of said load and weight values.

1 15. The server cluster of Claim 14 wherein each of said plurality of server  
2 computers include a policy data set store that provides for the storage of a distinct  
3 server configuration and wherein said load and weight values are dynamically  
4 determined by said plurality of server computers in response to said service  
5 requests based on said distinct server configurations of said plurality of server  
6 computers.

1 16. The server cluster of Claim 15 wherein said distinct server configurations  
2 include the distinct identities of said plurality of server computers.

1 17. The server cluster of Claim 16 wherein said distinct server configurations  
2 include distinct policy data relative to said service requests, wherein said host  
3 computer system is operative to collect, relative to respective said service requests,  
4 and provide attribute data to said plurality of server computers, and wherein said  
5 server computers evaluate said attribute data in conjunction with said distinct  
6 policy data to determine said weight values.

1 18. The server cluster of Claim 17 wherein said plurality of server computers  
2 implement a security processing service, wherein said host computer system is

3 operative to selectively route network transported data through said server  
4 computers dependent on said service requests as evaluated by said plurality of  
5 server computers.

1 19. The server cluster of Claim 18 said host computer is operative to initiate  
2 respective data transfer transactions for each of said service requests, wherein the  
3 default routing of each said data transfer transaction initially provides for the  
4 transfer of corresponding ones of said service requests to respective ones of said  
5 plurality of server computers, and wherein said respective ones of said plurality  
6 of server computers determine whether the subsequent routing of network data  
7 within said respective data transfer transactions includes routing said network data  
8 within said respective data transfer transactions through said plurality of server  
9 computers.

1 20. A computer system providing, on behalf of client computer systems, a  
2 network service through a scalable cluster of server computer systems, said system  
3 comprising:

4 a) a plurality of server computers coupled to provide a defined service,  
5 wherein a server computer of said plurality provides a response, including load  
6 information, in acknowledgment of a predetermined service request issued to said  
7 server computer system, said response selectively indicating nonacceptance of  
8 said predetermined service request; and

9 b) a client computer system having an identification list of said plurality of  
10 server computer systems, said client computer system being operative to  
11 autonomously select a first server computer system from said identification list to

12 which to issue said predetermined service request, wherein said client computer  
13 system is reactive to said response, on indicated nonacceptance of said  
14 predetermined service request, to autonomously select a second server computer  
15 system from said identification list to which to issue said predetermined service  
16 request, and wherein said client computer system is responsive to said load  
17 information of said response in subsequently autonomously selecting said first and  
18 second server computer systems.

1 21. The computer system of Claim 20 wherein said response further includes  
2 weight information and wherein said client computer system evaluates the  
3 combination of said load and weight information in autonomously selecting server  
4 computer systems from said identification list.

1 22. The computer system of Claim 21 wherein said plurality of server computer  
2 systems include respective policy engines and wherein said weight information  
3 reflects an association between a server computer policy role and said  
4 predetermined service request.

1 23. The computer system of Claim 22 wherein said predetermined service  
2 request includes predetermined client process attribute information and wherein  
3 said respective policy engines are responsive to said predetermined client process  
4 attribute information in determining said server computer policy role relative to  
5 said predetermined service request.

1 24. The computer system of Claim 23 wherein said load information includes  
2 a value representing network and server processor performance.

1 25. A method of dynamically managing the distribution of client requests to a  
2 plurality of server computer systems providing a network service, each of said  
3 server computer systems being discretely configured to respond to client requests,  
4 said method comprising the steps of:

5 a) processing client requests to select for a particular client request a  
6 particular server computer system of said plurality of server computer systems to  
7 service said particular client request, wherein the selection of said particular server  
8 computer system is dependent on the evaluation of accumulated selection  
9 qualification information;

10 b) forwarding said particular client request to said particular server  
11 computer system; and

12 c) receiving from said particular server computer system with respect to said  
13 particular client request instance selection qualification information discretely  
14 determined by said particular server computer system with respect to said  
15 particular client request, wherein said instance selection qualification information  
16 is incorporated into said accumulated selection qualification information.

1 26. The method of Claim 25 wherein said processing step dynamically  
2 evaluates said particular client request with respect to said accumulated selection  
3 qualification information to identify said particular server computer system as a  
4 best choice of said plurality of server computer systems for selection.



1 27. The method of Claim 26 further comprising the step of evaluating by said  
2 particular server computer system, subject to the discrete configuration of said  
3 particular server computer system, said particular client request to provide said  
4 instance selection qualification information.

1 28. The method of Claim 27 wherein said step of evaluating provides for the  
2 dynamic generation of said instance selection qualification information including  
3 a load value reflective of the performance capability of said particular server  
4 computer system.

1 29. The method of Claim 28 wherein said instance selection qualification  
2 information includes a relative prioritization of said particular client request with  
3 respect to said particular server computer system.

1 30. The method of Claim 29 wherein said client requests are issued with  
2 respect to client computer systems, wherein said particular client request includes  
3 attributes descriptive of a particular client computer system that issued said  
4 particular client request, and wherein said relative prioritization reflects the  
5 evaluation of said attributes with respect to said particular server computer system.

1 31. A method of distributing computational load over a plurality of server  
2 systems provided to support execution of a data processing service on behalf of  
3 a plurality of client systems, wherein the computational load is generated in  
4 response to client requests issued through a plurality of client processes, said  
5 method comprising the steps of:

6           a) first processing a particular client request to associate attribute data from  
7     a respective client process of said plurality of client processes with said particular  
8     client request;

9           b) selecting, for said particular client request, a particular target server  
10    system from among said plurality of server systems by matching said particular  
11    client request against accumulated selection information to identify said particular  
12    target server system;

13          c) second processing said particular client request, including said attribute  
14    data, by said particular target server system to dynamically generate instance  
15    selection information including a load value for said particular target server  
16    system and reflective of the combination of said particular client request and said  
17    particular target server system; and

18          d) incorporating said instance selection information into said accumulated  
19    selection information for subsequent use in said step of selecting.

1    32.    The method of Claim 31 wherein said instance selection information  
2    includes a relative weighting value reflective of the combination of said particular  
3    client request and said particular target server system and wherein said step of  
4    selecting matches said particular client request, including said attribute data,  
5    against corresponding data of said accumulated selection information to choose  
6    said particular target server system based on a best corresponding combination  
7    of relative weighting value and load value.

1    33.    The method of Claim 32 wherein said step of selecting includes a step of  
2    aging said accumulated selection information.

1 34. The method of Claim 33 further comprising the steps of:

2 a) first providing, through a host process, said particular client request,  
3 including attribute data, to said particular target server system; and

4 b) receiving by said host process, a particular target server response  
5 including said instance selection information;

6 c) determining, by said host process from said particular target server  
7 response, whether to select an alternate target server system;

8 d) reselecting, for said particular client request, a secondary target server  
9 system from among said plurality of server systems by matching said particular  
10 client request against said accumulated selection information, including said  
11 instance selection information received from said particular target server response  
12 to identify said secondary target server system; and

13 e) second providing, through said host process, said particular client  
14 request, including attribute data, to said alternate target server system.

1 35. The method of Claim 34 wherein said host process is executed on a client  
2 computer system.

1 36. The method of Claim 35 wherein said host process is executed on a  
2 gateway computer system coupleable through a communications network with a  
3 plurality of client computer systems.